

Dear Author,

Here are the proofs of your article.

- You can submit your corrections **online**, via **e-mail** or by **fax**.
- For **online** submission please insert your corrections in the online correction form. Always indicate the line number to which the correction refers.
- You can also insert your corrections in the proof PDF and **email** the annotated PDF.
- For fax submission, please ensure that your corrections are clearly legible. Use a fine black pen and write the correction in the margin, not too close to the edge of the page.
- Remember to note the **journal title**, **article number**, and **your name** when sending your response via e-mail or fax.
- **Check** the metadata sheet to make sure that the header information, especially author names and the corresponding affiliations are correctly shown.
- **Check** the questions that may have arisen during copy editing and insert your answers/ corrections.
- **Check** that the text is complete and that all figures, tables and their legends are included. Also check the accuracy of special characters, equations, and electronic supplementary material if applicable. If necessary refer to the *Edited manuscript*.
- The publication of inaccurate data such as dosages and units can have serious consequences. Please take particular care that all such details are correct.
- Please **do not** make changes that involve only matters of style. We have generally introduced forms that follow the journal's style. Substantial changes in content, e.g., new results, corrected values, title and authorship are not allowed without the approval of the responsible editor. In such a case, please contact the Editorial Office and return his/her consent together with the proof.
- If we do not receive your corrections **within 48 hours**, we will send you a reminder.
- Your article will be published **Online First** approximately one week after receipt of your corrected proofs. This is the **official first publication** citable with the DOI. **Further changes are, therefore, not possible.**
- The **printed version** will follow in a forthcoming issue.

Please note

After online publication, subscribers (personal/institutional) to this journal will have access to the complete article via the DOI using the URL: [http://dx.doi.org/\[DOI\]](http://dx.doi.org/[DOI]).

If you would like to know when your article has been published online, take advantage of our free alert service. For registration and further information go to: <http://www.springerlink.com>.

Due to the electronic nature of the procedure, the manuscript and the original figures will only be returned to you on special request. When you return your corrections, please inform us if you would like to have these documents returned.

Metadata of the article that will be visualized in OnlineFirst

ArticleTitle	TER-Plus: paraphrase, semantic, and alignment enhancements to Translation Edit Rate	
--------------	---	--

Article Sub-Title		
-------------------	--	--

Article CopyRight	Springer Science+Business Media B.V. (This will be the copyright line in the final PDF)	
-------------------	--	--

Journal Name	Machine Translation	
--------------	---------------------	--

Corresponding Author	Family Name	Snover
	Particle	
	Given Name	Matthew G.
	Suffix	
	Division	Laboratory for Computational Linguistics and Information Processing, Institute for Advanced Computer Studies
	Organization	University of Maryland
	Address	College Park, USA
	Email	snover@umiacs.umd.edu

Author	Family Name	Madnani
	Particle	
	Given Name	Nitin
	Suffix	
	Division	Laboratory for Computational Linguistics and Information Processing, Institute for Advanced Computer Studies
	Organization	University of Maryland
	Address	College Park, USA
	Email	nmadnani@umiacs.umd.edu

Author	Family Name	Dorr
	Particle	
	Given Name	Bonnie
	Suffix	
	Division	Laboratory for Computational Linguistics and Information Processing, Institute for Advanced Computer Studies
	Organization	University of Maryland
	Address	College Park, USA
	Email	bonnie@umiacs.umd.edu

Author	Family Name	Schwartz
	Particle	
	Given Name	Richard
	Suffix	
	Division	
	Organization	BBN Technologies
	Address	Cambridge, USA
	Email	schwartz@bbn.com

Schedule	Received	15 May 2009
	Revised	

Abstract

This paper describes a new evaluation metric, TER-Plus (TERP) for automatic evaluation of machine translation (MT). TERP is an extension of Translation Edit Rate (TER). It builds on the success of TER as an evaluation metric and alignment tool and addresses several of its weaknesses through the use of paraphrases, stemming, synonyms, as well as edit costs that can be automatically optimized to correlate better with various types of human judgments. We present a correlation study comparing TERP to BLEU, METEOR and TER, and illustrate that TERP can better evaluate translation adequacy.

Keywords (separated by '-')

Machine translation evaluation - Paraphrasing - Alignment

Footnote Information

Journal: 10590
Article: 9062



Author Query Form

**Please ensure you fill out your response to the queries raised below
and return this form along with your corrections**

Dear Author

During the process of typesetting your article, the following queries have arisen. Please check your typeset proof carefully against the queries listed below and mark the necessary changes either directly on the proof/online grid or in the 'Author's response' area provided below

Query	Details required	Author's response
1.	Please confirm the inserted city and country name are correct in Affiliation	
2.		

TER-Plus: paraphrase, semantic, and alignment enhancements to Translation Edit Rate

Matthew G. Snover · Nitin Madnani ·
Bonnie Dorr · Richard Schwartz

Received: 15 May 2009 / Accepted: 16 November 2009
© Springer Science+Business Media B.V. 2009

Abstract This paper describes a new evaluation metric, TER-Plus (TERP) for automatic evaluation of machine translation (MT). TERP is an extension of Translation Edit Rate (TER). It builds on the success of TER as an evaluation metric and alignment tool and addresses several of its weaknesses through the use of paraphrases, stemming, synonyms, as well as edit costs that can be automatically optimized to correlate better with various types of human judgments. We present a correlation study comparing TERP to BLEU, METEOR and TER, and illustrate that TERP can better evaluate translation adequacy.

Keywords Machine translation evaluation · Paraphrasing · Alignment

M. G. Snover (✉) · N. Madnani · B. Dorr
Laboratory for Computational Linguistics and Information Processing, Institute for Advanced
Computer Studies, University of Maryland, College Park, USA
e-mail: snover@umiacs.umd.edu

N. Madnani
e-mail: nmadnani@umiacs.umd.edu

B. Dorr
e-mail: bonnie@umiacs.umd.edu

R. Schwartz
BBN Technologies, Cambridge, USA
e-mail: schwartz@bbn.com

1 Introduction

TER-Plus, or TERP¹ (Snover et al. 2009), is an automatic evaluation metric for machine translation (MT) that scores a translation (the *hypothesis*) of a foreign language text (the *source*) against a translation of the source text that was created by a human translator, which we refer to as a *reference* translation. The set of possible correct translations is very large, possibly infinite, and any one reference translation represents a single point in that space. Frequently, multiple reference translations—typically 4—are provided to give broader sampling of the space of correct translations. Automatic MT evaluation metrics compare the hypothesis against this set of reference translations and assign a score to the similarity, such that a better score is given when the hypothesis is more similar to the references.

TERP follows this methodology and builds upon an already existing evaluation metric, Translation Error Rate (TER) (Snover et al. 2006). In addition to assigning a score to a hypothesis, TER provides an alignment between the hypothesis and the reference, enabling it to be useful beyond general translation evaluation. While TER has been shown to correlate well with translation quality, it has several flaws: it only considers exact matches when measuring the similarity of the hypothesis and the reference, and it can only compute this measure of similarity against a single reference. The handicap of using a single reference can be addressed by constructing a lattice of reference translations—this technique has been used to combine the output of multiple translation systems (Rosti et al. 2007). TERP does not utilize this methodology and instead addresses the exact matching flaw of TER.

In addition to aligning words in the hypothesis and reference if they are exact matches, TERP uses stemming and synonymy to allow matches between words. It also uses probabilistic phrasal substitutions to align phrases in the hypothesis and reference. These phrase substitutions are generated by considering possible paraphrases of the reference words. Matching using stems and synonyms (Banerjee and Lavie 2005) as well as using paraphrases (Zhou et al. 2006; Kauchak and Barzilay 2006) have been shown to be beneficial for automatic MT evaluation. Paraphrases have been shown to be additionally useful in expanding the number of references used for evaluation (Madnani et al. 2008) although they are not used in this fashion within TERP. The use of synonymy, stemming, and paraphrases allows TERP to better cope with the limited number of reference translations provided. TERP was one of the top metrics submitted to the NIST Metrics MATR 2008 challenge (Przybocki et al. 2008), having the highest average rank over all the test conditions (Snover et al. 2009).

We first discuss the original TER metric in Sect. 2. In Sect. 3, we present the details of our various enhancements to TER. We then briefly review the alignment capability of TERP along with some examples in Sect. 4. Finally, in Sect. 5, we show the results of optimizing TERP for human judgments of adequacy and compare with other established evaluation metrics, followed by an analysis of the relative benefits of each of the new features of TERP in Sect. 6.

¹ TERP is named after the nickname—“terp”—of the University of Maryland, College Park, Mascot: the diamondback terrapin.

51 2 Translation Edit Rate (TER)

52 One of the first automatic metrics used to evaluate automatic MT systems was Word
 53 Error Rate (WER) (Niessen et al. 2000), which remains the standard evaluation met-
 54 ric for Automatic Speech Recognition. WER is computed as the Levenshtein distance
 55 between the words of the hypothesis and the words of the reference translation divided
 56 by the length of the reference translation. Unlike speech recognition, there are many
 57 correct translations for any given foreign sentence. These correct translations differ
 58 not only in lexical choice but also in the order in which the words occur. WER is
 59 inadequate for evaluating MT output as it fails to combine knowledge from multiple
 60 reference translations and cannot accurately model the reordering of words and phrases
 61 that accompanies translation.

62 TER addresses the latter failing of WER by allowing block movement of words,
 63 called *shifts*, within the hypothesis. Shifting a phrase is assumed to have the same *edit*
 64 *cost* as inserting, deleting or substituting a word, regardless of the number of words
 65 being shifted. While a general solution to WER with block movements is NP-Com-
 66 plete (Lopresti and Tomkins 1997), TER computes an approximate solution by using
 67 a greedy search to select the words to be shifted, as well as imposing additional con-
 68 straints on these words. These constraints are intended to simulate the way in which
 69 a human editor might choose the words to shift. Other automatic metrics exist that
 70 have the same general formulation as TER but address the complexity of shifting in
 71 different ways, such as the CDER evaluation metric (Leusch et al. 2006).

72 The shifting constraints used by TER serve to better model the quality of translation
 73 as well as to reduce the model's computational complexity. Examining a larger set of
 74 shifts, or choosing them in a more optimal fashion might result in a lower TER score
 75 but it would not necessarily improve the ability of the measure to determine the quality
 76 of a translation. The constraints used by TER are as follows:

- 77 1. Shifts are selected by a greedy algorithm that chooses the shift that yields the
- 78 largest reduction in WER between the reference and the hypothesis.
- 79 2. The sequence of words shifted in the hypothesis must *exactly match* the sequence
- 80 of words in the reference that it will align with after the shift.
- 81 3. The words being shifted, and the matching reference words, must each contain
- 82 at least one error, according to WER, before the shift occurs. This prevents the
- 83 shifting of words that are already correctly matched.

84 When TER is used with multiple references, it does not combine the references, but,
 85 instead, scores the hypothesis against each reference individually—as is the case with
 86 metrics such as METEOR (Banerjee and Lavie 2005). The reference against which
 87 the hypothesis has the fewest number of edits is deemed to be the closest reference, and
 88 the final TER score is the number of edits in between the hypothesis and this closest
 89 reference divided by the average number words across all of the references.

90 3 TER-Plus (TERp)

91 TER-Plus extends the TER metric beyond the limitation of exact matches through the
 92 addition of three new types of edit operations, detailed in Sect. 3.1: stem matches,

93 synonym matches, and phrase substitutions using automatically generated para-
 94 phrases. These changes allow a relaxing of the shifting constraints used in TER,
 95 which is discussed in Sect. 3.2. In addition, instead of all edit operations having a
 96 uniform edit cost of 1—as is the case in TER—the edit costs for TERP can be learned
 97 automatically in order to maximize correlation with human judgments. The details of
 98 this optimization are presented in Sect. 3.3.

99 3.1 Stem, synonym, and paraphrase substitutions

100 In addition to the edit operations of TER—Matches, Insertions, Deletions, Substitu-
 101 tions and Shifts—TERP also uses three new edit operations: Stem Matches, Synonym
 102 Matches and Phrase Substitutions. Rather than treating all substitution operations as
 103 edits of cost 1, the cost of a substitution in TERP varies so that a lower cost is used if
 104 two words are synonyms (a Synonym Match), share the same stem (a Stem Match), or
 105 if two phrases are paraphrases of each other (a Phrase Substitution). The cost of these
 106 new edit types is set, along with the other edit costs, according to the type of human
 107 judgment for which TERP is optimized, as described in Sect. 3.3.

108 TERP identifies stems and synonyms in the same manner as the METEOR met-
 109 ric (Banerjee and Lavie 2005), where words are determined to share the same stem
 110 using the Porter stemming algorithm (Porter 1980), and words are determined to be
 111 synonyms if they share the same synonym set according to WordNet (Fellbaum 1998).

112 Phrase substitutions are identified by looking up—in a pre-computed *phrase*
 113 *table*—probabilistic paraphrases of phrases in the reference to phrases in the hypoth-
 114 esis. The paraphrases used in TERP are automatically extracted using the pivot-
 115 based method (Bannard and Callison-Burch 2005) with several additional filtering
 116 mechanisms to increase precision. The pivot-based method identifies paraphrases as
 117 English phrases that translate to the same foreign phrase in a bi-lingual phrase table.
 118 The corpus used for paraphrase extraction was an Arabic-English newswire bi-text
 119 containing a million sentences, resulting in a phrase table containing approximately
 120 15 million paraphrase pairs. While an Arabic-English corpus was used to generate the
 121 paraphrases, the resulting phrase pairs are English only and can be applied to regard-
 122 less of the source language. We have previously shown that the choice of data for
 123 paraphrasing is not of vital importance to TERP's performance (Snover et al. 2009).
 124 A few examples of the extracted paraphrase pairs that were actually used by TERP in
 125 experiments described later are shown below:

brief \Rightarrow short controversy over \Rightarrow polemic about by using power \Rightarrow by force response \Rightarrow reaction

126
 127 Some paraphrases, such as *brief* and *short* are redundant with other edit types used by
 128 TERP such as synonym and stem matching.

129 A probability for each paraphrase pair is estimated as described in Bannard and
 130 Callison-Burch (2005). However, studies (Snover et al. 2009) of these paraphrase
 131 probabilities have shown that they are not always reliable indicators of the semantic
 132 relatedness of phrase pairs and further refinements of these probability estimates might
 133 prove valuable to TERP and other MT evaluation metrics.

134 With the exception of the phrase substitutions, all of the edit operations used by
 135 TERP have fixed cost edits, i.e., the edit cost is the same regardless of the words
 136 in question. The cost of a phrase substitution is a function of the probability of the
 137 paraphrase and the number of edits needed to align the two phrases without the use of
 138 phrase substitutions. In effect, the probability of the paraphrase is used to determine
 139 how much to discount the alignment of the two phrases. For a phrasal substitution
 140 between a reference phrase r and a hypothesis phrase h where Pr is the probability of
 141 paraphrasing r as h , and $\text{edit}(r, h)$ is number of edits needed to align r and h without
 142 any phrasal substitutions, the edit cost is specified by three parameters, w_1 , w_2 , and
 143 w_3 as follows:

$$144 \quad \text{cost}(r, h) = w_1 + \text{edit}(r, h)(w_2 \log(\text{Pr}) + w_3)$$

145 Only paraphrases specified in the input phrase table are considered for phrase sub-
 146 stitutions. In addition, the total cost for a phrasal substitution is limited to values
 147 greater than or equal to 0, to ensure that the edit cost for substitution operations is
 148 always non-negative. The parameter w_1 allows a constant cost to be specified for all
 149 phrase substitutions, while parameters w_2 and w_3 adjust the discount applied to the
 150 edit cost of the two phrases.

151 3.2 Additional differences from TER

152 In addition to the new edit operations, TERP differs from TER in several other ways.
 153 First, TERP is insensitive to casing information since we observe that penalizing
 154 for errors in capitalization lowers the correlation with human judgments of translation
 155 quality. Second, TERP is capped at 1.0. While the formula for TER allows it to exceed
 156 1.0 if the number of edits exceed the number of words, such a score would be unfair
 157 since the hypothesis cannot be more than 100% wrong.

158 The shifting criteria in TERP have also been relaxed relative to TER, so that shifts
 159 are allowed if the words being shifted are: (i) exactly the same, (ii) synonyms, stems
 160 or paraphrases of the corresponding reference words, or (iii) any such combination. In
 161 addition, a set of stop words is used to constrain the shift operations such that common
 162 words (“the”, “a” etc.) and punctuation can be shifted if and only if a non-stop word is
 163 also shifted. This reduces the number of shifts considered in the search and prevents
 164 any shifts that may not correspond with an increase in translation quality.

165 More relaxed shift constraints have been explored that allowed shifts even if some
 166 words did not match at all. We have empirically found this greatly increased the
 167 number of shifts considered, but also significantly decreased correlation with human
 168 judgment. The shift constraints imposed by TER and TERP serve not only to speed up

169 the algorithm but also correspond to those block movement of words that correspond
170 with increased translation quality.

171 3.3 TERP edit cost optimization

172 While TER uses uniform edit costs—1 for all edits except matches—, we seek to
173 improve TERP’s correlation with human judgments by weighting different edit types
174 more heavily than others, as some types of errors are more harmful to translation
175 quality than others.

176 TERP uses a total of eight edit costs. However, the cost of an exact match is held
177 fixed at 0 which leaves a total of seven edit costs that can be optimized. Since the para-
178 phrase edit cost is represented by three parameters, this yields a total of nine parameters
179 that are varied during optimization. All parameters, except for the three phrasal substi-
180 tution parameters, are also restricted to be positive. A hill-climbing search optimizes
181 the parameters to maximize the correlation of human judgments with the TERP score.
182 In this paper, these correlations are measured at the sentence, or *segment*, level. How-
183 ever, optimization could also be performed to maximize document level correlation
184 or any other measure of correlation with human judgments.

185 While TERP can be run using a fixed set of parameters, it can be beneficial to tune
186 them depending on the properties of translation desired. Optimization of MT evalua-
187 tion metrics towards specific human judgment types has previously investigated in a
188 similar manner by Lita et al. (2005). Depending on whether the end goal is to maxi-
189 mize correlation with HTER, adequacy, or fluency, different sets of parameters may
190 better reflect translation performance (Snover et al. 2009).

191 4 TERp alignment

192 In addition to providing a score indicating the quality of a translation, TERP gener-
193 ates an alignment between the hypothesis and reference, indicating which words are
194 correct, incorrect, misplaced, or similar to the reference translation. While the quality
195 of this alignment is limited by the similarity of the reference to the hypothesis it can
196 be beneficial in diagnosing error types in MT systems.

197 Actual examples of TERP alignments are shown in Fig. 1. Within each example,
198 the first line is the reference translation, the second line is the original hypothesis,
199 and the third line is the hypothesis after performing all shifts. Words in **bold** are
200 shifted, while square brackets are used to indicate other edit types: *P* for phrase sub-
201 stitutions, *T* for stem matches, *Y* for synonym matches, *D* for deletions, and *I* for
202 insertions.

203 These alignments allow TERP to provide quality judgments on translations and
204 to serve as a diagnostic tool for evaluating particular types of translation errors.
205 In addition, it may also be used as a general-purpose string alignment tool—TER
206 has been used for aligning multiple system outputs to each other for MT system
207 combination (Rosti et al. 2007), a task for which TERP may be even better
208 suited.

R : ... [a number of]_D leaders expressed their opposition to [participating in]_P the government ...
 H : ... the leaders expressed their opposition to the government **take part in** ...
 H' : ... [the]_I leaders expressed their opposition to [**take part in**]_P the government ...

5.27 (6)

R : ... [he]_D [went on to say]_P , "we also discussed how [to galvanize]_D the ... "
 H : ... continued , "we also discussed how the activation of ... "
 H' : ... [continued]_P , "we also discussed how the [activation of]_I ... "

6.48 (8)

R : ... [but]_{S1} we [have]_{Y1} [Palestinian]_T [,]_{S2} Arab [or]_D Islamic [alternatives]_{Y2} .
 H : ... and we now possess an **Islamic** or the Palestinians and Arab options.
 H' : ...[and]_{S1} we now [possess]_{Y1} [an or the]_I [Palestinians]_T [and]_{S2} Arab **Islamic** [options]_{Y2} .

6.14 (10)

Fig. 1 Examples of TERP alignment output. In each example, **R**, **H** and **H'** denote the reference, the original hypothesis and the hypothesis after shifting respectively. Shifted words are **bolded** and other edits are in *[brackets]*. Number of edits shown: TERP (TER)

209 5 Experimental results

210 5.1 Optimization for adequacy

211 In order to tune and test TERP, we used a portion of the Open MT-Eval 2006 evalu-
 212 ation set that had been annotated for adequacy (on a seven-point scale) and released
 213 by NIST as a development set for the Metrics MATR 2008 challenge (Przybocki et al.
 214 2008). This set consists of the translation hypotheses from 8 Arabic-to-English MT
 215 systems for 25 documents, which in total consisted of 249 segments. For each segment,
 216 four reference translations were also provided. Optimization was done using 2-fold
 217 cross-validation. These optimized edit costs (and subsequent results) differ slightly
 218 from the formulation of TERP submitted to the Metrics MATR 2008 challenge, where
 219 tuning was done without cross-validation. Optimization requires small amounts of
 220 data but should be done rarely so that the metric can be held constant to aid in system
 221 development and comparison.

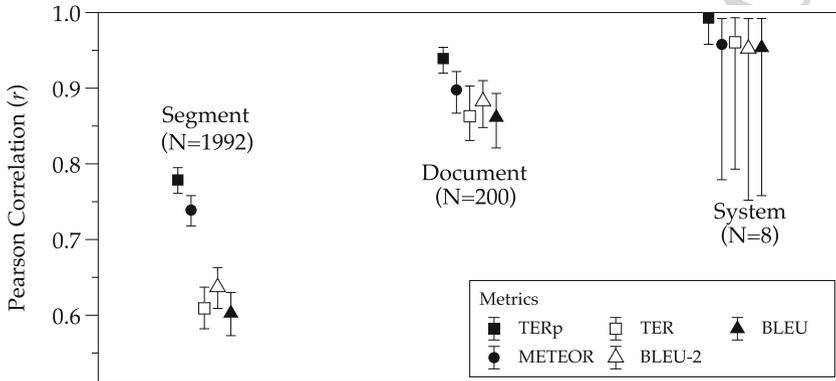
222 TERP parameters were then optimized to maximize segment level Pearson corre-
 223 lation with adequacy on the tuning set. The optimized edit costs, averaged between
 224 the two splits of the data, are shown in Table 1. Because segment level correlation
 225 places equal importance on all segments, this optimization over-tunes for short seg-
 226 ments, as they have very minor effect at the document or system level. Optimization
 227 on length weighted segment level correlation would rectify this but would result in
 228 slightly worse segment level correlations.

229 5.2 Correlation results

230 In our experiments, we compared TERP with METEOR (Banerjee and Lavie 2005)
 231 (version 0.6 using the Exact, WordNet synonym, and Porter stemming modules), TER

Table 1 TERP edit costs optimized for adequacy

Match	Insert	Deletion	Substitution	Stem
0.0	0.20	0.97	1.04	0.10
Syn.	Shift	Phrase substitution		
0.10	0.27	$w_1: 0.0$	$w_2: -0.12$	$w_3: 0.19$

**Fig. 2** Metric correlations with adequacy on the Metrics MATR 2008 development set. Correlations are significantly different if the center point of one correlation does not lie within the confidence interval of the other correlation

(version 0.7.25), the IBM version of BLEU (Papineni et al. 2002) with a maximum n -gram size of 4 (BLEU). We also included a better correlating variant of BLEU with a maximum n -gram size of 2 (BLEU-2). TER and both versions of BLEU were run in case insensitive mode as this produces significantly higher correlations with human judgments, while METEOR is already case insensitive.

To evaluate the quality of an automatic metric, we examined the Pearson correlation of the automatic metric scores—at the segment, document and system level—with the human judgments of adequacy. Document and system level adequacy scores were calculated using the length weighted averages of the appropriate segment level scores.

Pearson correlation results between the automatic metrics and human judgments of adequacy are shown in Fig. 2. We can determine whether the difference between two correlation coefficients is statistically significant by examining the confidence interval of the Pearson coefficient, r . If the correlation coefficient for a metric occurs within the 95% confidence interval of another metric, then the difference between the correlations of the metrics is not statistically significant.

TERP consistently outperformed all of the other metrics on the segment, document, and system level Pearson correlations, with all but one difference being statistically significant. While TERP had higher correlation than TER on the system level, the

² Confidence intervals are calculated using the Fisher r -to- z transformation, consulting a z -table to find the upper and lower bounds of a 95% confidence interval, and then converting the values back to r scores. This is solely a function of the correlation coefficient, r , and the number of data points, N .

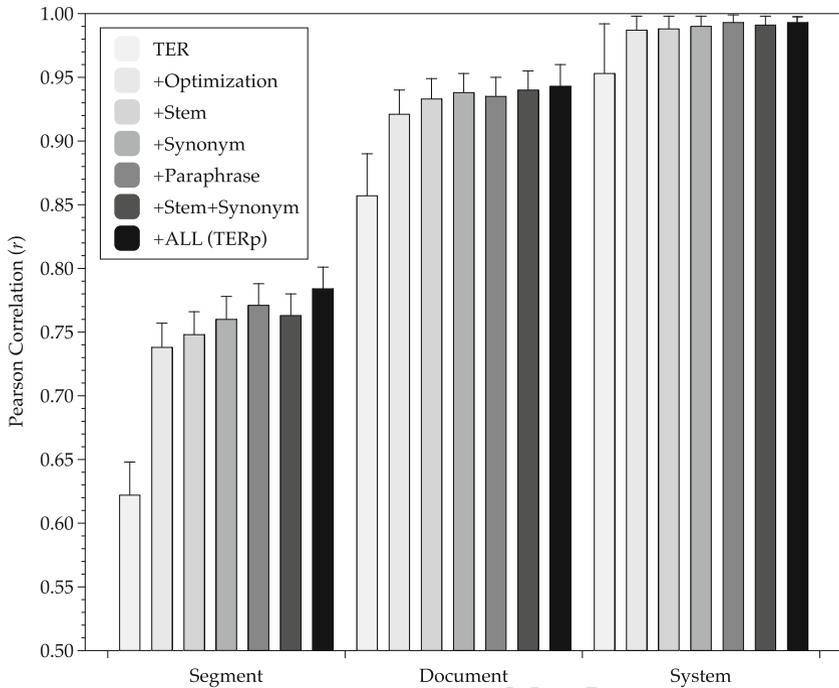


Fig. 3 Pearson correlation of TERP with selective features

250 difference is not statistically significant—the differences with all other metrics are
 251 statistically significant. Of the other metrics, METEOR consistently had the highest
 252 Pearson correlation at the segment and document level. METEOR, the only other tun-
 253 able metric, might possibly correlate better by retuning for this dataset, although this
 254 is not generally done for METEOR.

255 6 Benefit of individual TERp features

256 In this section, we examine the benefit of each of the new features of TERP by indi-
 257 vidualy adding each feature to TER and measuring the correlation with the human
 258 judgments. Each condition was optimized as described in Sect. 5.1. Figure 3 shows the
 259 Pearson correlations for each experimental condition along with the 95% confidence
 260 intervals.

261 The largest gain over TER is through the addition of optimizable edit costs. This
 262 takes TER from being a metric with balanced insertion and deletion costs to a recall-
 263 oriented metric which strongly penalizes deletion errors, while being forgiving of
 264 insertion errors. This single addition gives statistically significant improvements over
 265 TER at the segment and document levels. This validates similar observations of the
 266 importance of recall noted by Lavie et al. (2004).

267 The other three features of TERP—stemming, synonymy, and paraphrases—are
 268 added on top of the optimized TER condition since optimization is required to

269 determine the edit costs for the new features. The addition of each of these features
 270 increases correlations over the optimized edit costs at all levels, with statistically sig-
 271 nificant gains at the segment level for the addition of synonymy or paraphrasing. The
 272 addition of paraphrasing gives the largest overall gains in correlation after optimization
 273 and is more beneficial than stemming and synonymy combined. A large percentage of
 274 synonym and stem matches are already captured in the paraphrase set and, therefore,
 275 the combination of all three features yields only a small gain over paraphrasing alone.

276 The TERP framework and software also provides for separate word classes with
 277 individual edit costs, so that the edit costs of various sets of words can be increased
 278 or decreased. For example, the cost of deleting content words could be set higher than
 279 that of deleting function words. It is difficult to set such costs manually as it is not
 280 clear how these phenomenon are treated by human annotators of translation quality,
 281 although these costs could be determined by automatic optimization.

282 7 Discussion

283 TERP extends the TER metric using stems, synonyms, paraphrases, and optimizable
 284 edit costs to assign a more realistic score to a translation hypothesis and generate a
 285 better alignment against the reference. Experimental results show that TERP achieves
 286 significant gains in correlation with human judgments over other MT evaluation met-
 287 rics. Optimization can be targeted towards specific types of human judgments, yielding
 288 different edit costs for TERP, for use in cases when a specific notion of translation
 289 quality is desired.

290 Automatic MT evaluation metrics are used for two major purposes: (1) To compare
 291 two or more MT systems (or variants of the same system) in order to determine which
 292 system generates better translations. This is often used to show that the addition of a
 293 new feature to a translation system yields an improvement over a baseline system. (2)
 294 To automatically optimize or tune the parameters of a system. While we conducted
 295 this study in the context of the first purpose—showing that TERP provides significant
 296 gains in evaluating final system outputs—we have not evaluated TERP for the second
 297 purpose. It is frequently the case that automatic metrics that appear useful according
 298 to the first criterion are not suitable for the second purpose, resulting in degenerate
 299 system parameters. To evaluate a metric's suitability for optimization, a translation
 300 system must be optimized using a baseline metric, such as BLEU, and also using the
 301 new metric being examined. The final outputs of the two systems tuned to the different
 302 metrics must then be judged by humans to determine which optimization method pro-
 303 vides better translations. Unfortunately, this technique is also biased by the translation
 304 system that is being tuned and the method used for parameter optimization. Further
 305 explorations of this nature are needed to determine if TERP, and other metrics, are
 306 suitable for use in MT parameter optimization.

307 We showed that the addition of stemming, synonymy and, most importantly,
 308 paraphrasing to the TER metric significantly improves its correlation with human
 309 judgments. We believe that further significant improvements in TERP and other auto-
 310 matic evaluation metrics are contingent on the use of additional linguistic features so

311 as to better capture the fluency of a translation hypothesis and its similarity in meaning
312 to reference translations.

313 **Acknowledgments** This work was supported, in part, by BBN Technologies under the GALE Program,
314 DARPA/IPTO Contract No. HR0011-06-C-0022 and in part by the Human Language Technology Center of
315 Excellence. The authors would like to thank Philip Resnik, Chris Callison-Burch, Mark Przybocki, Sebas-
316 tian Bronsart and Audrey Le. The TERP software is available online for download at: [http://www.umiacs.
umd.edu/~snover/terp/](http://www.umiacs.
317 umd.edu/~snover/terp/).

318 References

- 319 Banerjee S, Lavie A (2005) METEOR: an automatic metric for MT evaluation with improved correla-
320 tion with human judgments. In: Proceedings of the ACL 2005 workshop on intrinsic and extrinsic
321 evaluation measures for MT and/or summarization, pp 228–231
- 322 Bannard C, Callison-Burch C (2005) Paraphrasing with bilingual parallel corpora. In: Proceedings of
323 the 43rd annual meeting of the association for computational linguistics (ACL 2005). Ann Arbor,
324 Michigan, pp 597–604
- 325 Fellbaum C (1998) WordNet: an electronic lexical database. MIT Press. [http://www.cogsci.princeton.edu/
wn](http://www.cogsci.princeton.edu/
326 wn). Accessed 7 Sep 2000
- 327 Kauchak D, Barzilay R (2006) Paraphrasing for automatic evaluation. In: Proceedings of the human
328 language technology conference of the North American chapter of the ACL, pp 455–462
- 329 Lavie A, Sagae K, Jayaraman S (2004) The significance of recall in automatic metrics for MT evaluation.
330 In: Proceedings of the 6th conference of the association for machine translation in the Americas,
331 pp 134–143
- 332 Leusch G, Ueffing N, Ney H (2006) CDER: efficient MT evaluation using block movements. In: Proceed-
333 ings of the 11th conference of the European chapter of the association for computational linguistics,
334 pp 241–248
- 335 Lita LV, Rogati M, Lavie A (2005) BLANC: learning evaluation metrics for MT. In: Proceedings of human
336 language technology conference and conference on empirical methods in natural language processing
337 (HLT/EMNLP). Vancouver, BC, pp 740–747
- 338 Lopresti D, Tomkins A (1997) Block edit models for approximate string matching. *Theor Comput Sci*
339 181(1):159–179
- 340 Madnani N, Resnik P, Dorr BJ, Schwartz R (2008) Are multiple reference translations necessary? Investi-
341 gating the value of paraphrased reference translations in parameter optimization. In: Proceedings of
342 the eighth conference of the association for machine translation in the Americas, pp 143–152
- 343 Niessen S, Och F, Leusch G, Ney H (2000) An evaluation tool for machine translation: fast evaluation for MT
344 research. In: Proceedings of the 2nd international conference on language resources and evaluation,
345 pp 39–45
- 346 Papineni K, Roukos S, Ward T, Zhu W-J (2002) Bleu: a method for automatic evaluation of machine trans-
347 lation. In: Proceedings of the 40th annual meeting of the association for computational linguistics,
348 pp 311–318
- 349 Porter MF (1980) An algorithm for suffix stripping. *Program* 14(3):130–137
- 350 Przybocki M, Peterson K, Bronsart S (2008) Official results of the NIST 2008 “Metrics for MACHine
351 TRanslation” Challenge (MetricsMATR08). <http://nist.gov/speech/tests/metricsmatr/2008/results/>
- 352 Rosti A-V, Matsoukas S, Schwartz R (2007) Improved word-level system combination for machine trans-
353 lation. In: Proceedings of the 45th annual meeting of the association of computational linguistics.
354 Prague, Czech Republic, pp 312–319
- 355 Snover M, Dorr B, Schwartz R, Micciulla L, Makhoul J (2006) A study of translation edit rate with targeted
356 human annotation. In: Proceedings of association for machine translation in the Americas, pp 223–231
- 357 Snover M, Madnani N, Dorr B, Schwartz R (2009) Fluency, adequacy, or HTER? Exploring different human
358 judgments with a tunable MT metric. In: Proceedings of the fourth workshop on statistical machine
359 translation. Association for Computational Linguistics, Athens, Greece, pp 259–268
- 360 Zhou L, Lin C-Y, Hovy E (2006) Re-evaluating machine translation results with paraphrase support. In:
361 Proceedings of the 2006 conference on empirical methods in natural language processing (EMNLP
362 2006), pp 77–84